



Pretests for genetic-programming evolved trading programs: “zero-intelligence” strategies and lottery trading

Shu-Heng Chen, Nicolas Navet

► To cite this version:

Shu-Heng Chen, Nicolas Navet. Pretests for genetic-programming evolved trading programs: “zero-intelligence” strategies and lottery trading. The 13th International Conference on Neural Information Processing - ICONIP2006, Oct 2006, Hong-Kong, Hong Kong SAR China. pp.450-460, 10.1007/11893295_50 . inria-00105926

HAL Id: inria-00105926

<https://hal.inria.fr/inria-00105926>

Submitted on 27 Aug 2007

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L’archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d’enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Pretests for genetic-programming evolved trading programs: “zero-intelligence” strategies and lottery trading^{*}

Shu-Heng Chen¹ and Nicolas Navet^{1,2}

¹ AI-ECON Research Center, Department of Economics,
National Chengchi University, Taipei, Taiwan 11623,
`chchen@nccu.edu.tw` ,

² LORIA-INRIA, Campus-Scientifique, BP239, F-54506 Vandoeuvre, France,
`nnavet@loria.fr`

Abstract. Over the last decade, numerous papers have investigated the use of GP for creating financial trading strategies. Typically in the literature results are inconclusive but the investigators always suggest the possibility of further improvements, leaving the conclusion regarding the effectiveness of GP undecided. In this paper, we discuss a series of pretests, based on several variants of random search, aiming at giving more clear-cut answers as to whether a GP scheme, or any other machine-learning technique, can be effective with the training data at hand. Precisely, pretesting allows us to distinguish between a failure due to the market being efficient or due to GP being inefficient. The analysis is illustrated with GP-evolved strategies for three stock exchanges exhibiting different trends.

1 Motivation and Introduction

The computational intelligence techniques such as genetic programming³, with their continuous advancement, persistently bring us something positive to expect, and incessantly push the application domain to more challenging issues. However, sometimes, the costs and benefits of using these advanced CI techniques are uncertain. Usually the benefits are not assured, while the costs are immediate. On the one hand, the CI techniques are frequently used as intensive search algorithms, which are inevitably computationally demanding, and take up a great amount of computational resources. On the other hand, whether there is a needle in the haystack remains dubious.⁴ Certainly, if such a needle does not

^{*} An extended version of this paper will appear as chapter 8 of the book *Computational Intelligence in Economics and Finance - Volume 2* , Springer Verlag, to be published in 2007.

³ Although, in this paper, we only focus on genetic programming, but the general ideas and some specific implementations may also be applicable to other computational intelligence techniques used to induce trading strategies.

⁴ For example, in the financial application domain, the lack of such a needle may be due to the efficient market hypothesis or the no-arbitrage condition.

exist at all, the all efforts are made to no avail. Given this asymmetry between costs and benefits, it would be economical, at the first stage, to test the existence of such a needle before a fully-fledged version of search is applied. We refer to this procedure as a *pretest*.

The pretest procedure proposed here is in a sense similar to the pretests used in econometrics where the estimator of an unknown parameter is chosen on the basis of the outcome of a pretest ([1]). Pretesting, also known as “data-snooping” in finance, classically serves to select the right model that will be used later on for forecasting purpose ([2,3]). More broadly, pretesting can be considered to be a practice of a sequential decision-making process, which is used when the decision involves a great deal of uncertainty, and the costs of making a wrong decision are huge.⁵ In this case, at the first stage, we would like to expend some limited resources in probing into gaining some initial information, e.g. the distribution of a very uncertain environment, while in the later stages, we will make our decision based on the gauged distribution.

The reasoning behind pretesting is very intuitive, and [4] is the first to apply this idea to the financial application of genetic programming (GP). [4] proposed a measure known as the η statistic. The η statistic is a measure of predictability. Basically, using a simple (vanilla) version of GP, one can first gauge the predictability based on η . When η is low or close to zero, it indicates that there is nothing to forecast. So, the use of fully-fledged GP is not advised. The virtue of doing this is to distinguish *two kinds of possibilities* when we see a failure of an initial attempt based on simple GP. First, the series itself has nothing to forecast; second, GP has not been used appropriately. Understanding this distinction can result in big differences in our second stage of the decision. In the former case, we may simply give up any further search to avoid wasting resources. In the latter case, we should keep on exploring different deliberations of GP to search for potential gains before a final conclusion can be made. In either case, we have a clear-cut situation. However, when a pretest is absent, we become less conclusive: we are no longer sure whether the problem is due to the non-existence of the needle, or the improper use of GP.

Unfortunately, in most financial trading applications of GP, a pretest has been largely neglected.⁶ We think that this negligence may give rise to many inconclusive results. Typically, what happens is that the results from using GP are not very convincing, but the investigators always suggest directions for further

⁵ The problem of sequential decision making under incomplete knowledge has been studied by researchers in various fields, such as optimal control, psychology, economics, and game theory.

⁶ This may not be completely so. In fact, most earlier studies selected the buy-and-hold strategies or a risk-free investment (e.g., treasury bills) as the benchmark. However, the conclusion that “GP performs better than buy-and-hold in a bearish market and worse in a bullish market” is often found in the literature. However, nothing different can be expected since buy-and-hold is the worst possible strategy in a steadily decreasing market and the best possible strategy in a steadily increasing market. This shows the limits of choosing buy-and- hold as a benchmark. See, for example, [5].

improvement, leaving the actual conclusion regarding the effectiveness of GP undecided. Therefore, this study attempts to provide practical pretesting procedure aimed at reducing the number of cases where the conclusion is inconclusive.

Needless to say, there are various ways of implementing different types of pretesting. For example, the η statistic mentioned above can be used as a pretest, as [4] did, but that is mainly applied to forecasting time series. That a series is to a certain extent predictable does not necessarily imply that we can develop profitable trading strategies. For example, the predictability horizon might be too short, the fluctuation might not be volatile enough to cover the round-trip transaction costs or, simply, the right trading instrument might not be available (e.g., no short selling in a downward oriented market) or else they are some regulation and rules (e.g., the “uptick rule” makes intraday trading with short selling more difficult). Consequently, literature on forecasting with GP and literature on trading with GP usually are separated. Therefore, in this paper, we attempt to develop pretest procedures that are more suitable for trading purposes.

More precisely, we will propose several different styles of pretests, which when put together can help us decide whether there are hidden patterns to be discovered and whether GP is properly designed to do the job. The essential idea underlying all proposed pretests is to compare the performance of GP with random trading strategies or behavior. However, as we shall see in Section 2, just making trading strategies or trading behavior arbitrarily random is not sufficient to provide a fair and informative comparison. To do so, some constraints are expected, and the intriguing point is how to impose these constraints properly.

The rest of the paper is organized as follows. Section 2 provides a detail formulation of four pretests. The first three are concerned with the trading strategies, whereas the last one is concerned with the trading behavior. Normally, trading behavior comes from trading strategies, and they cannot be separated; however, when randomness is introduced, difference between the two may arise. In particular, in the vein of algorithmic complexity, random trading strategies can imply trading behavior actually using knowledge, while random trading behavior presumably excludes such a possibility. We, therefore, intentionally distinguish between the two by calling the former *zero-intelligence strategies*, and the latter *lottery trading*. Section 3 discuss how to use these proposed tests together to make a better judgement given the initial results we have. Section 4 illustrates the proposed pretests based on the real detail and the experimental designs detailed in the appendix. Section 5 gives the concluding remarks.

2 Pretests: description and rationale

In this section, we describe a series of 4 pretests and discuss their purpose and implementation. Of the 4 pretests, we highlight that 2 are of particular interest and, as shown in Section 3, enable us to gain complementary knowledge on the data under study and on the efficiency of the GP’s implementation. In the following, we consider GP with a validation stage before the actual testing on the

out-of-sample data. Validation means that the best rules induced on the training interval are further selected on the unseen data, the validation period, before being applied out-of-sample. The validation step is a device to fight overfitting⁷ that has been widely used in earlier GP work (see for instance [7,8]). Note that our pretest proposals remain valid for GP without the validation step except that pretest 2 replaces pretest 1, which requires validation.

2.1 GP versus equivalent intensity random search

The basic idea here is to compare the outcome of GP with an *equivalent intensity random search*. We say that two search algorithms are equivalent in terms of search intensity if their execution leads to the evaluation of the same number of different trading strategies on the training data. For instance, let us consider GP with the parameters chosen for this study: a population of 500 individuals evolved over 100 generations. In the first approximation, the equivalent random search (ERS) would consist of evaluating 50,000 randomly created solutions. In practice, search algorithms sometimes rediscover identical solutions over the course of their execution. This can be detected by keeping track of all created individuals since the beginning of the execution, and, in doing so, useless fitness evaluations can be skipped, which actually saves computing time when the fitness function is rather time-consuming as it is in our context. Since, computationally speaking, what is preponderant is the fitness evaluation, and since the extent to which GP re-discovers the same individuals is very dependent upon the implementation, we impose that our definition of equivalent search intensity only accounts for unique individuals, *i.e.* individuals which require evaluation. We consider two solutions to be different if their expression is *syntactically different*⁸, in our context, if the trees representing the programs are different.

The three following pretests compare GP with a random search both with and without training and validation stage. In the latter search technique, the biologically inspired evolution process of GP is simply replaced by the creation of solutions at random. Since with random search the strategies do not benefit

⁷ The actual effectiveness of validation in this context is, however, still an open question, see [5] and [6].

⁸ Two individuals can be syntactically different while being equivalent in the sense that they lead to equivalent trading decisions, the equivalence could thus also be defined in terms of semantics. With symbolic simplification using rewriting rules and interval arithmetic on the function arguments, we could detect that some syntactically different individuals are in fact semantically identical. However, there is no way of making sure that all duplicates will be detected and the implementation of this procedure would be so complex and time consuming at run time that, in our opinion, a definition based on semantics would be of little practical interest. Alternatively, the equivalence in search intensity could be defined in terms of equivalent computing time. However there is such a difference of complexity between a fully-fledged GP implementation and random search that it is hard to imagine how we can ensure that the two implementations have been optimized in a similar manner, while a better implementation of GP for instance may lead us to come to an opposite conclusion.

from the “intelligence” resulting from the evolution or learning process, we dub randomly created solutions as *zero-intelligence trading strategies*.

For each pretest i , we formulate the null hypothesis $\mathcal{H}_{i,0}$ that GP does not outperform the technique it is compared with at pretest i , where the alternative hypothesis is denoted by $\mathcal{H}_{i,1}$. The experiments will provide us with the answer on whether $\mathcal{H}_{i,0}$ should be rejected in favour of $\mathcal{H}_{i,1}$ or not.

Pretest 1: GP versus equal search intensity random search with training and validation stage. The implementation of the random search strategy is straightforward: parameters of GP are set in such a way that only the initial generation, where individuals are created at random, is used. The size of the initial population is adjusted so that the resulting search intensity is identical to the one of the regular GP.

- **Hypothesis $\mathcal{H}_{1,0}$ cannot be rejected:** the first explanation that can be envisaged is that, GP or not, there is nothing essential to be learned from the past. In that case GP would strongly “overfit” the training data, possibly explaining that its out-of-sample performance is worse than with a random search. This can be due by the market being efficient or because the training interval is very dissimilar to the out-of-sample⁹. Another explanation is that the GP machinery is not working properly, for instance due to a wrong choice of the function/terminal sets, because the parameters are inappropriate (e.g. too low search intensity), or the genetic operators are unable to create better-than-random individuals.
- **Hypothesis $\mathcal{H}_{1,0}$ is rejected in favour of $\mathcal{H}_{1,1}$:** there may be something to learn from the past and GP, with the chosen parameters, may be effective in that task.

Rejecting $\mathcal{H}_{1,0}$ is of course a first indication of the efficiency of GP but we cannot rule out the case where there would be nothing useful to learn on the data at hand and GP would beat random search by mere luck. We will see in Section 3, that further investigation may provide additional evidence to answer that question.

Pretest 2: GP versus equal search intensity random search with training but without a validation stage. Here, the best random solutions on the training interval are applied directly to the out-of-sample period. With regard to pretest 1, pretest 2 could give us some insight about how effective is validation as a device to fight against overfitting. However, since overfitting is unlikely to occur with random solutions, the rationale for using pretest 2 is unclear and it will not be further considered in this study. A more direct and effective way to evaluate the effect of the validation stage is simply to compare regular GP with and without validation¹⁰.

⁹ In [5], numerous experiments have highlighted that when training and out-of-sample data sets are very “dissimilar”, for instance if the market exhibits an opposite trend, then there is little chance that GP will perform well out-of-sample.

¹⁰ For instance, as it is done in [5].

Pretest 3: GP versus equal search intensity random search without training and without validation stage. In pretest 3, the selection of the strategies on the training set is removed: a large number random strategies are created and applied directly out-of-sample. The performance is evaluated as the average performance (e.g. average total return) over the set of random strategies. Comparing the outcome of pretest 3 with regard to pretest 1 and regular GP tells us something about how effective the selection process is, the extent to which a top performing rule on the training and validation sets will keep on performing well out-of-sample. If strategies selected by GP or random search on the training and validation intervals have some predictive ability out-of-sample, it provides use with some evidence that there is something to learn from the past. It is worth pointing out that the randomness of the strategies is here constrained by the GP language: rules can only be made with GP functions/terminals organized according to the typing scheme. For instance, it is possible that the GP language is not sufficiently expressive to define a rule consisting in buying and selling every other period¹¹. In the remainder of this study, we will consider pretest 4, presented in Section 2.2, that is similar in spirit to pretest 3, but is more random in the sense that it does not possess the bias in randomness induced by the GP language.

2.2 GP versus lottery trading

We refer to *lottery trading* as a strategy that would consist in making the investment decision at random on the basis of the outcome of a random variable. In its simplest form, the random variable would follow a Bernoulli distribution where the parameter p expresses the probability to take a long position and $1 - p$ the probability to be out of the market.

In our context, this requires refinement since we are interested in profitability and profitability takes into account transaction costs. So, to allow a fair comparison with GP, we should make sure that the expected number of transactions for lottery trading is the same as for GP. We refer to the expected number of transactions per unit of time as the *frequency of a trading strategy*. Another important characteristic of a trading strategy is what we term its *intensity*, *i.e.* the number of periods where a position¹² “in the market” is held, over the length of the trading interval. We should also enforce lottery trading to have the same expected intensity as GP to avoid misleading results, for instance, the case where, given its frequency, the intensity of lottery trading is not sufficient to cover the transaction costs with the volatility of the market under study.

¹¹ Period refers to the granularity of time used for trading, for instance, one second or one day.

¹² Implicitly, we consider here the trading of a single instrument, e.g. an index, where 2 decisions are possible at each time period: be in or be out of the market without short selling, or with short selling as implemented in [5], holding a long position or a short position. The concept remains valid where one can be holding a long position, a short position or be out of the market. One can also define the intensity and the frequency of a strategy for each instrument traded.

One denotes by F_{GP} and I_{GP} respectively the average frequency and average intensity observed for the set of GP evolved rules applied on the testing interval over all GP runs, N_{GP} is the number of transactions leading to F_{GP} . For the experiments made in the following, a sequence of investment decisions S_{LT} resulting from lottery trading is generated at random according to the following procedure:

- the intensity for lottery trading, I_{LT} , is uniformly chosen in $[I_{GP} \cdot (1 - \alpha), \min(1, I_{GP} \cdot (1 + \alpha))]$ with $0 \leq \alpha \leq 1$. In a first step, S_{LT} is made of the '0' positions (*i.e.* out of the market) followed by the block of '1' positions (*i.e.* in the market) corresponding to I_{LT} ,
- the number of transactions N_{LT} is uniformly chosen in the set of integer values that are even¹³ in interval $[N_{GP} \cdot (1 - \alpha), N_{GP} \cdot (1 + \alpha)]$. The block of '1' is subdivided at random in $N_{LT}/2$ sub-sequences and each sub-sequence is inserted at random inside the block of '0'. This design avoids the problem of overlapping among the '1' sub-sequences that may occur with other schemes.

We formulate the pretest comparing GP and lottery trading and denote by $\mathcal{H}_{4,0}$ the null hypothesis that GP does not outperform lottery trading.

Pretest 4: GP versus lottery trading. Obviously, if GP is not able to outperform lottery trading, it gives strong evidence that GP will not be good at evolving effective trading strategies with the data at hand. In Section 3, we shall discuss this point in more detail.

3 What does the pretest tell us ?

The outcomes of the pretests provide us with answers to the following two questions: is there something essential to learn on the training data that can be of interest for the out-of-sample period ? Does the GP implementation show some evidence of effectiveness in that task ? Clearly, before actually trading with GP evolved programs, these two questions must be answered with reasonable certainty; the rest of this section explains how pretests may help in that regard.

3.1 Question 1: is there something to learn ?

The null hypothesis $\mathcal{H}_{4,0}$ corresponding to pretest 4 has been presented in Section 2.2. We introduce pretest 5 that will be used in conjunction with pretest 4.

¹³ N_{LT} has to be even since a "buy" transaction is followed by a sell transaction and no positions are left open.

Pretest 5: equivalent intensity random search with training and validation versus lottery trading. Here, we compare lottery trading to a random search with training and validation, and a search intensity equivalent to the one used for GP in pretest 4. The null hypothesis $\mathcal{H}_{5,0}$ is that the equivalent intensity random does not outperform lottery trading on the out-of-sample data. Depending on the validity of $\mathcal{H}_{4,0}$ and $\mathcal{H}_{5,0}$, we can draw the conclusions that are summarized in Table 1.

	$\mathcal{H}_{4,0}$	$\mathcal{H}_{5,0}$	Interpretation
case 1	$\neg\mathbf{R}$	$\neg\mathbf{R}$	there is evidence that there is nothing to learn
case 2	\mathbf{R}	$\neg\mathbf{R}$	there may be something to learn (weak certainty)
case 3	\mathbf{R}	\mathbf{R}	there is evidence that there is something to learn
case 4	$\neg\mathbf{R}$	\mathbf{R}	there may be something to learn (weak certainty)

Table 1. Information drawn from the outcomes of pretest 4 and pretest 5 ($\neg\mathbf{R}$ means that the null hypothesis $\mathcal{H}_{i,0}$ cannot be rejected while \mathbf{R} means that the hypothesis is rejected in favour of the alternative hypothesis).

In case 1, best solutions on the training intervals, obtained with 2 different search algorithms, do not perform better than lottery trading on the out-of-sample period. This suggests to us that there is nothing to learn. In case 2, GP outperforms lottery trading but random search does not; it is possible that there is something to learn, but that the selected random rules do not have a sufficient predictive ability. Anyway, this leads us to a less certain conclusion than in case 3 where both search techniques outperform lottery trading. Finally case 4 is a special case where random search performs better than lottery trading but GP does not. The whole evolutionary process of GP has thus a detrimental effect and a possible explanation is that GP-induced solutions overfit the training data.

3.2 Question 2: is the GP machinery working properly ?

The second question we ought to ask is whether GP is effective. Of course, this cannot be answered with the data at hand if pretests 4 and 5 have shown that there is nothing to be learned (case 1 in Table 1). In addition, in case 4 of Table 1, we already know that GP is not efficient since, by transitivity, it is outperformed by the random search-based algorithm. Thus, the only two cases where one really needs to proceed to further examination are case 2 and case 3. The validity of the null hypothesis $\mathcal{H}_{1,0}$, which can be tested with pretest 1, gives a helpful insight into the answer: only if $\mathcal{H}_{1,0}$ should be rejected can we conclude that GP shows some real effectiveness. We would like to stress that rejecting $\mathcal{H}_{1,0}$ is far from implying profitability, but beating a mere random search algorithm on a difficult problem with an infinite search space is the bare minimum one can expect from GP.

4 Experiments

The aim of the experiments is to evaluate the extent to which the pretests proposed are reliable. The methodology adopted here is to check if the outcomes of the pretests are consistent with results already published in the literature. We call GP2 the GP implementation developed for this study and GP1 the software¹⁴ used in [5], which will constitute our benchmark. The GP2 control parameters, as close as possible to the ones used in [5] for GP1, are summarized in Table 1 (Appendix A).

The traded instruments are the indexes of 3 stock exchanges: the TSE 300 (Canada), the Nikkei Dow Jones (Japan) and the Capitalization Weighted Stock Index (Taiwan). They have been chosen among the 8 markets studied in [5] because they exhibit the main evolution patterns that can be found in the set of 8 markets. The aim of GP is to induce the most profitable strategy, measured by the accumulated return, for trading the stock exchange index. The use of short selling is possible. We adopt what is done classically in literature in terms of data-preprocessing and use normalized data that is obtained by dividing each day's price by a 250-day moving average¹⁵. In a way similar to what is done usually, we subdivide the whole dataset into three sections: the *training*, *validation* and *out-of-sample* test periods. For each stock index considered, 3 different out-of-sample test periods of 2 years each (*i.e.* 1999-2000, 2001-2002, 2003-2004) follow a 3-year validation and a 3-year training period. In the following, the term market refers to a stock exchange during a specific out-of-sample period. For instance, market Canada-1 (C1 for short) is the TSE 300 during the out-of-sample period 1999-2000. Hypothesis testing is performed with the *Student's t-test* at a 95% confidence level. The samples for statistics are constituted of the results of 50 GP runs, 50 runs of equivalent search intensity random search with training and validation (ERS) and 100 runs of lottery trading (LT).

In 4 out of the 9 markets (*i.e.* C3, J2, T1, T3), there is evidence that there is something to learn from the training data (case 3 in Table 1). This is consistent with [5] where GP2 performs outstandingly on these 4 markets (respective total return: 0.34, 0.17, 0.52, 0.27 with GP1). In markets C1, J3 and T2, pretests 4 and 5 suggest to us that there is nothing to learn (case 1). Except for C1, GP2 also performs poorly (-0.18 for J3 and -0.05 for T2). Finally, in the 3 markets where GP2 is shown to beat ERS ($\mathcal{H}_{1,0}$ is rejected in favor of $\mathcal{H}_{1,1}$ for J1, J2 and T1), the GP results are very good : both GP1 and GP2 produce positive returns and outperform the buy-and-hold strategy.

Although more comprehensive tests are needed, the experiments conducted here show some preliminary evidence that the proposed pretests possess some predictive ability. Indeed, when the outcome is “nothing to learn”, the two GP

¹⁴ Although both programs have been developed by members of the AI-ECON Research Center, they have not been written by the same persons and do not share a single line of code. Furthermore GP2, which is based on the *Open-Beagle* library (see <http://beagle.gel.ulaval.ca/>), implements strongly-typed GP.

¹⁵ See [5] for a discussion about how non-normalized data affects the performance of GP.

implementation perform very poorly (except in one case). On the contrary, when the pretests suggest that there is something to learn, at least one implementation does well, and when GP2 is more efficient than random search (*i.e.* ERS), GP1 from [5] is efficient too. In the light of the pretests, we should also conclude that our GP implementation (*i.e.* GP2) is more efficient than ERS (GP2 outperforms ERS in 3 markets while ERS never beats GP2 with statistical significance). However, in our experiments, searching trading rules at random, with the same set of functions and terminals as used in GP, is usually enough to come up with trading systems that outperform lottery trading when GP does as well. This suggests to us that GP2 may only be able to take advantage of “simple” regularities in the data.

5 Conclusions

The main purpose of this paper is to enrich the earlier research on Genetic Programming (GP) induced market-timing decisions by proposing pretests aiming to shed light on the GP results. In actual fact, in the literature, the results of applying GP for market-timing decisions are typically not very convincing, but the investigators always suggest the possibility of further improvements. If the investigators can first convince that there is something to learn and that GP is suitable for that task, then their conclusion would be less vague and uncertain. We propose here a series of pretests, where GP is tested against a random behavior (*lottery trading*) and against strategies created at random (*zero-intelligence strategies*), that aim to answer these two crucial questions. Of course there is the risk of getting a wrong pretest result and the possible reasons why GP may have failed should be thoroughly investigated before drawing a conclusion. But, in the end, analyzing the results in the light of the pretests should help draw more fine-grained conclusions.

Acknowledgment. This research was conducted when the second author (Nicolas Navet) was visiting researcher at the AI-ECON Research Center, National Chengchi University (NCCU), Taipei, Taiwan. The financial support from the AI-ECON Research Center as well as NCCU and INRIA is greatly acknowledged. The authors would like also to acknowledge the grant from National Science Council #95-2415-H-004-002-MY3.

References

1. J.A. Giles and D.E.A. Giles: Pre-test Estimation and Testing in Econometrics: Recent Developments. *Journal of Economic Surveys* **7**(2) (1993) 145–97
2. D. Danilov and J.R. Magnus: Forecast Accuracy After Pretesting with an Application to the Stock Market. *Journal of Forecasting* **23** (2004) 251–274
3. R. Sullivan and A. Timmermann and H. White: Data-Snooping, Technical Trading Rule Performance, and the Bootstrap. *Journal of Finance* **54** (1999) 1647–1692
4. M.A. Kaboudan: A Measure of Time Series Predictability Using Genetic Programming Applied to Stock Returns. *Journal of Forecasting* **18** (1999) 345–357

5. S.-H. Chen and T.-W. Kuo and K.-M. Hoi: Genetic Programming and Financial Trading: How Much about "What we Know". In: Handbook of Financial Engineering. Kluwer Academic Publishers (2006) Forthcoming.
6. S.-H. Chen and T.-Z. Kuo: Overfitting or Poor Learning: A Critique of Current Financial Applications of GP. In Springer-Verlag, ed.: Proceedings of the Sixth European Conference on Genetic Programming (EuroGP-2003). Number 2610 in LNCS (2003) 34–46
7. F. Allen and R. Karjalainen: Using Genetic Algorithms to Find Technical trading rules. Journal of Financial Economics **51** (1999) 245–271
8. C. Neely and P. Weller and R. Dittmar: Is Technical Analysis in the Foreign Exchange Market Profitable? A Genetic Programming Approach. Journal of Financial and Quantitative Analysis **32**(4) (1997) 405–427

A Genetic programming settings

Program GP2 implements strongly typed GP with the set of functions and terminals described in Table 1. The parameters here are basically identical to the ones in [5] (program GP1) except when fine-tuning GP2 have highlighted that better results may be obtained with different parameters. Precisely, we make use of more elitism, the size of the tournament selection is set to 5 and numerical mutation is implemented.

Population size	500
Number of generations	100
Maxim tree depth	10
Function set	+, -, *, /, norm, average, max, min, lag, and, or, not, >, <, if-then-else, true, false
Terminal set	price, real and integer ephemeral constants
Value range for real constants	[-1,1]
Value range for integer constants	[0,1000]
Offsprings created by:	
crossover	50%
standard mutation	20%
swap mutation	15%
reproduction	10%
ephemeral constant mutation	5%
Initialization	ramp-half-and-half
Evolution scheme	generation-by-generation replacement strategy
Elitism	25 best individuals kept for next generation
Selection scheme	tournament selection of size 5
Fitness function	accumulated return
Transaction costs	0.5%
Validation	
number of best trees saved	1 individual per run is saved for validation

Fig. 1. GP control parameters